

日 本 国 特 許 庁  
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日

Date of Application:

2002年 7月19日

出 願 番 号

Application Number:

特願2002-211634

[ST.10/C]:

[JP2002-211634]

出 願 人

Applicant(s):

インターナショナル・ビジネス・マシーンズ・コーポレーション

2002年12月13日

特 許 庁 長 官  
Commissioner,  
Japan Patent Office

太田信一郎



出証番号 出証特2002-3098275

【書類名】 特許願

【整理番号】 JP9020109

【提出日】 平成14年 7月19日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 17/30

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内

【氏名】 野美山 浩

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内

【氏名】 岩男 俊孝

【特許出願人】

【識別番号】 390009531

【氏名又は名称】 インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】 100086243

【弁理士】

【氏名又は名称】 坂口 博

【代理人】

【識別番号】 100091568

【弁理士】

【氏名又は名称】 市位 嘉宏

【代理人】

【識別番号】 100108501

【弁理士】

【氏名又は名称】 上野 剛史

【復代理人】

【識別番号】 100104880

【弁理士】

【氏名又は名称】 古部 次郎

【手数料の表示】

【予納台帳番号】 081504

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9706050

【包括委任状番号】 9704733

【包括委任状番号】 0004480

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 情報検索システム、情報検索方法、HTML文書の構造解析方法及びプログラム

【特許請求の範囲】

【請求項1】 ネットワークを介してウェブサイトのクローリングを行う情報検索システムにおいて、

所定のウェブページにおける意味を考慮してソースコードの構造を解析する構造解析部と、

前記構造解析部の解析結果に基づいて、前記所定のウェブページからリンクされる他のウェブサイトの重要度を計算する重要度計算部と、

前記重要度計算部により計算された重要度に応じてウェブサイトをクロールするクローリング実行部と  
を備えることを特徴とする情報検索システム。

【請求項2】 前記構造解析部は、前記ソースコードに含まれる情報要素のうち、相互に関連する情報要素を対応付けることを特徴とする請求項1に記載の情報検索システム。

【請求項3】 前記重要度計算部は、前記ウェブサイトの重要度を計算するための戦略を、予め用意された戦略の中から選択的に用いて重要度の計算を行うことを特徴とする請求項1に記載の情報検索システム。

【請求項4】 前記重要度計算部は、前記ウェブサイトの重要度を計算するための戦略として、複数の戦略を選択し、各々重みを付けて用いることを特徴とする請求項3に記載の情報検索システム。

【請求項5】 HTML文書の文書構造を意味を考慮して解析し、解析によって得られた情報要素を対応するアンカーに付加する文書構造解析部と、

前記文書構造解析部の解析により得られた前記情報要素に基づいて計算された前記アンカーの重要度に応じて当該アンカーにてリンクされるウェブサイトをクロールするクローリング実行部と  
を備えることを特徴とする情報検索システム。

【請求項6】 前記文書構造解析部は、前記HTML文書を構成する各情報

要素を、当該情報要素が持つ意味のまとまりごとにブロック化し、各ブロック内の情報要素を付加情報として同一ブロック内のアンカーに付加することを特徴とする請求項 5 に記載の情報検索システム。

【請求項 7】 前記文書構造解析部の解析により得られた前記情報要素に基づき、予め選択された所定の戦略にしたがって前記アンカーの重要度を計算する重要度計算部をさらに備え、

前記クロール実行部は、前記重要度計算部にて計算された前記アンカーの重要度に応じてウェブサイトをクロールすることを特徴とする請求項 5 に記載の情報検索システム。

【請求項 8】 コンピュータを用いて、ネットワークを介してウェブサイトのクロールを行う情報検索方法であって、

初期情報となるウェブページを取得してソースコードを記憶装置に格納するステップと、

前記記憶装置から前記ウェブページのソースコードを読み出し、当該ウェブページにおける意味を考慮して構造解析を行い、解析結果を前記記憶装置に格納するステップと、

前記記憶装置に格納された前記構造解析の結果に基づいて、前記ウェブページからリンクされる他のウェブサイトの重要度を計算するステップと、

計算された重要度に応じてウェブサイトにアクセスし、コンテンツを取得して前記記憶装置に格納するステップと  
を含むことを特徴とする情報検索方法。

【請求項 9】 メモリから処理対象である HTML 文書を読み出し、当該 HTML 文書を構成する各情報要素を、当該 HTML 文書のタグに基づいてブロック化し、ブロック化された当該 HTML 文書の構造データをメモリに格納する第 1 のステップと、

前記メモリからブロック化された前記 HTML 文書の構造データを読み出し、意味的に相互に関連する情報要素を対応付けることにより、当該 HTML 文書のブロック構造を更新し、更新された当該構造データをメモリに格納する第 2 のステップと

を含むことを特徴とするコンピュータを用いたHTML文書の構造解析方法。

【請求項10】 ネットワークに接続されたコンピュータを制御して、ウェブサイトのクローリングを行うプログラムであって、

初期情報となるウェブページを取得してソースコードを記憶装置に格納する処理と、

前記記憶装置から前記ウェブページのソースコードを読み出し、当該ウェブページにおける意味を考慮して構造解析を行い、解析結果を前記記憶装置に格納する処理と、

前記記憶装置に格納された前記構造解析の結果に基づいて、前記ウェブページからリンクされる他のウェブサイトの重要度を計算する処理と、

計算された重要度に応じてウェブサイトにアクセスし、コンテンツを取得して前記記憶装置に格納する処理と

を前記コンピュータに実行させることを特徴とするプログラム。

【請求項11】 前記プログラムは、前記ソースコードに含まれる情報要素のうち、相互に関連する情報要素を対応付けることにより、前記構造解析を前記コンピュータに実行させることを特徴とする請求項10に記載のプログラム。

【請求項12】 前記プログラムによる前記ウェブサイトの重要度を計算する処理では、前記ウェブサイトの重要度を計算するための戦略として、複数の戦略を各々重みを付けて用いることを特徴とする請求項10に記載のプログラム。

【請求項13】 コンピュータを制御して、HTML文書の文書構造を解析するプログラムであって、

メモリから処理対象であるHTML文書を読み出し、当該HTML文書を構成する各情報要素を、当該HTML文書のタグに基づいてブロック化し、ブロック化された当該HTML文書の構造データをメモリに格納する第1の処理と、

前記メモリからブロック化された前記HTML文書の構造データを読み出し、意味的に相互に関連する情報要素を対応付けることにより、当該HTML文書のブロック構造を更新し、更新された当該構造データをメモリに格納する第2の処理と

を前記コンピュータに実行させることを特徴とするプログラム。

【請求項 1 4】 前記プログラムによる前記第 2 の処理では、  
文書構造解析の目的に鑑みて不要な情報要素を識別する処理と、  
構造的に意味のないブロックを削除する処理と、  
前記情報要素の内容に基づいて、情報要素のマージあるいはブロックの分割を行う処理と、

各ブロックに含まれる情報に基づいて、ブロック構造をマージする処理と  
を前記コンピュータに実行させることを特徴とする請求項 1 3 に記載のプログラム。

【請求項 1 5】 ネットワークに接続されたコンピュータを制御してウェブサイトのクローリングを行うプログラムを、当該コンピュータが読み取り可能に記録した記録媒体において、

前記プログラムは、  
初期情報となるウェブページを取得してソースコードを記憶装置に格納する処理と、

前記記憶装置から前記ウェブページのソースコードを読み出し、当該ウェブページにおける意味を考慮して構造解析を行い、解析結果を前記記憶装置に格納する処理と、

前記記憶装置に格納された前記構造解析の結果に基づいて、前記ウェブページからリンクされる他のウェブサイトの重要度を計算する処理と、

計算された重要度に応じてウェブサイトにアクセスし、コンテンツを取得して前記記憶装置に格納する処理と  
を前記コンピュータに実行させることを特徴とする記録媒体。

【請求項 1 6】 コンピュータを制御して HTML 文書の文書構造を解析するプログラムを、当該コンピュータが読み取り可能に記録した記録媒体において、

前記プログラムは、  
メモリから処理対象である HTML 文書を読み出し、当該 HTML 文書を構成する各情報要素を、当該 HTML 文書のタグに基づいてブロック化し、ブロック化された当該 HTML 文書の構造データをメモリに格納する第 1 の処理と、

前記メモリからブロック化された前記HTML文書の構造データを読み出し、意味的に相互に関連する情報要素を対応付けることにより、当該HTML文書のブロック構造を更新し、更新された当該構造データをメモリに格納する第2の処理と

を前記コンピュータに実行させることを特徴とする記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、ネットワークを介して所望の情報を自動的に取得する技術に関し、特にインターネット上で提供されているウェブコンテンツを検索し、リンクを辿りながら取得（クローリング）する技術に関する。

【0002】

【従来の技術】

今日、インターネットに代表されるコンピュータのネットワーク環境が広く普及したことにより、ネットワーク上で提供されている膨大な情報の中から検索エンジンを用いて所望の情報を検索し取得することが一般的に行われている。この検索エンジンには多くの種類が存在するが、予め情報を検索して取得しておき、検索要求に応じて保持している情報を返す静的な検索エンジンを用いた場合は、膨大な情報源（ウェブページ等）を対象としなければならないため、最新の情報を獲得することが困難である。また、基本的に検索エンジンを持つサーバが全ての処理を行うことが前提となっているため、サーバの負担が大きい。

【0003】

そこで、静的な検索エンジンで集めたキーワード検索結果の集合を初期集合として用い、これを起点に関連するサイトを動的に検索する手法が提案されている。この種の従来技術としては、例えば、下記文献1に開示されたShark-Searchと呼ばれる検索技術がある。

文献1：Michael Herscovici, Michal Jacovi, Yoelle S. Maarek, Dan Pelleg, Menachem Shtalhaim, Sigalit Ur. "The Shark-Search Algorithm: An Applica



tion: Tailored Web SiteMapping” In the Proceedings of WWW7, the 7th International World Wide Web Conference, Brisbane, April 1998. Also appeared in the Journal of Computer Networks and ISDN 30 (1998), pp 317-326.

HYPERLINK ”<http://www7.scu.edu.au/programme/fullpapers/1849/com1849.htm>  
” <http://www7.scu.edu.au/programme/fullpapers/1849/com1849.htm>

【 0 0 0 4 】

同文献に開示された技術は、インターネット上で、指定されたURL (Uniform Resource Locator) とキーワードとに基づき、当該URLのウェブサイトから当該キーワードに関連するウェブサイト（重要度の高いウェブサイト）を動的に検索する。このシステムは、初期集合を求めるためのキーワード (domain query) と動的にウェブサイトをクロールする際にクロール対象であるウェブサイトの重要度の計算に用いるキーワード (focused query) の2つを用いることによって精度の向上を図っている。

【 0 0 0 5 】

【発明が解決しようとする課題】

上述したように、ネットワーク上で提供される膨大な最新の情報を効率よく検索するために、検索要求があった場合に動的に情報検索を行うことが求められる。

。

しかし、上述した従来の動的な検索エンジンは、基本的に、ユーザが指定したトピック（キーワード等）に近い情報という1つの判断基準 (relevance()) で検索を行う。そのため、情報の使用目的に応じて多様な戦略で柔軟に検索を行うことができなかった。

【 0 0 0 6 】

また、情報を効率的に検索するためには、取得対象である情報（ウェブページ等）の重要度を判断し、これに基づいて情報の取得順や取得範囲を決める必要がある。しかし、インターネット上でURL及びトピックに基づいてウェブサイトをクロールする従来技術では、この重要度を効果的に判断することができなかった。すなわち、情報の重要度を判断するために限られた情報、例えばウェブペー

ジにおける、指定されたキーワードや、アンカーに近い位置に記載されているテキストなどしか用いておらず、所望の情報を効率的に検索することができなかった。例えば、上記文献 1 に開示された従来技術の場合、文献 1 には、アンカーの重要性を判断するために当該アンカーの近傍のテキスト（`anchor_text_context`）を考慮するとの記述があるが、どのようにしてこの `anchor_text_context` を得るかについて明確な記述はない。

#### 【0007】

そこで、本発明は、情報の使用目的に応じて多様な戦略による柔軟な情報検索を可能とすることを目的とする。

また、本発明は、ウェブサイトのクローリングにおいて、この多様な戦略による情報検索を実現するために、ウェブページに含まれる情報を有効に活用して検索を行うことを目的とする。

#### 【0008】

##### 【課題を解決するための手段】

上記の目的を達成する本発明は、ネットワークを介してウェブサイトのクローリングを行う、次のように構成された情報検索システムとして実現される。すなわち、この情報検索システムは、所定のウェブページにおける意味を考慮してソースコードの構造を解析する構造解析部と、この構造解析部の解析結果に基づいて、このウェブページからリンクされる他のウェブサイトの重要度を計算する重要度計算部と、この重要度計算部により計算された重要度に応じてウェブサイトをクロールするクローリング実行部とを備えることを特徴とする。

より詳しくは、この重要度計算部は、ウェブサイトの重要度を計算するための戦略を、予め用意された戦略の中から選択的に用いて重要度の計算を行う。さらに好ましくは、この重要度計算部は、複数の戦略を選択し、各々重みを付けて用いる。

#### 【0009】

また、本発明による他の情報検索システムは、HTML 文書の文書構造を意味を考慮して解析し、解析によって得られた情報要素を対応するアンカーに付加する文書構造解析部と、この文書構造解析部の解析により得られた情報要素に基づ

いて計算されたアンカーの重要度に応じて、このアンカーにてリンクされるウェブサイトをクロールするクローリング実行部とを備えることを特徴とする。

ここで詳細には、この文書構造解析部は、HTML文書を構成する各情報要素を、この情報要素が持つ意味のまとまりごとにブロック化し、各ブロック内の情報要素を付加情報として同一ブロック内のアンカーに付加する。

この情報検索システムは、文書構造解析部の解析により得られた情報要素に基づき、予め選択された所定の戦略にしたがってアンカーの重要度を計算する重要度計算部をさらに備えることができる。これにより、クローリング実行部は、重要度計算部にて所定の戦略にしたがって計算されたアンカーの重要度に応じてウェブサイトをクロールすることができる。

#### 【0010】

また、上記の目的を達成する本発明は、コンピュータを用いネットワークを介してウェブサイトのクローリングを行う、次のような情報検索方法として実現される。この情報検索方法は、初期情報となるウェブページを取得してソースコードを記憶装置に格納するステップと、この記憶装置からウェブページのソースコードを読み出し、このウェブページにおける意味を考慮して構造解析を行うステップと、この構造解析の結果に基づいて、このウェブページからリンクされる他のウェブサイトの重要度を計算するステップと、計算された重要度に応じてウェブサイトにアクセスし、コンテンツを取得するステップとを含むことを特徴とする。

#### 【0011】

さらに、本発明は、この情報検索方法などで用いられる、次のようなHTML文書の構造解析方法としても実現される。すなわち、このHTML文書の構造解析方法は、処理対象であるHTML文書を構成する各情報要素を、このHTML文書のタグに基づいてブロック化するステップと、ブロック化された前記HTML文書の構造データにおいて、意味的に相互に関連する情報要素を対応付けることにより、このHTML文書のブロック構造を更新するステップとを含む。

より詳細には、このHTML文書のブロック構造を更新するステップは、文書構造解析の目的に鑑みて不要な情報要素を識別するステップと、構造的に意味の

ないブロックを削除するステップと、この情報要素の内容に基づいて、情報要素のマージあるいはブロックの分割を行うステップと、各ブロックに含まれる情報に基づいて、ブロック構造をマージするステップとを含む。

【 0 0 1 2 】

また、本発明は、コンピュータを制御して上述した情報検索システムとして機能させるプログラムや、上述した情報検索方法またはHTML文書の構造解析方法における各ステップに対応する処理をコンピュータに実行させるプログラムとして実現することができる。このプログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記録媒体に格納して配布したり、ネットワークを介して配信したりすることにより、提供することができる。

【 0 0 1 3 】

【発明の実施の形態】

以下、添付図面に示す実施の形態に基づいて、この発明を詳細に説明する。

図1は、本実施の形態による情報検索システムを実現するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

図1に示すコンピュータ装置は、演算手段であるCPU (Central Processing Unit: 中央処理装置) 101と、M/B (マザーボード) チップセット102及びCPUバスを介してCPU101に接続されたメインメモリ103と、同じくM/Bチップセット102及びAGP (Accelerated Graphics Port) を介してCPU101に接続されたビデオカード104と、PCI (Peripheral Component Interconnect) バスを介してM/Bチップセット102に接続されたハードディスク105、ネットワークインターフェイス106及びUSBポート107と、さらにこのPCIバスからブリッジ回路108及びISA (Industry Standard Architecture) バスなどの低速なバスを介してM/Bチップセット102に接続されたフロッピーディスクドライブ109及びキーボード/マウス110とを備える。

なお、図1は本実施の形態を実現するコンピュータ装置のハードウェア構成を例示するに過ぎず、本実施の形態を適用可能であれば、他の種々の構成を取ることができる。例えば、ビデオカード104を設ける代わりに、ビデオメモリのみ

を搭載し、CPU 101にてイメージデータを処理する構成としても良いし、ATA (AT Attachment) などのインターフェイスを介してCD-ROM (Compact Disc Read Only Memory) やDVD-ROM (Digital Versatile Disc Read Only Memory) のドライブを設けても良い。

#### 【0014】

本実施の形態では、情報としてインターネット上で提供される各種のウェブコンテンツ（ウェブページやそのオブジェクト）を検索し獲得する場合を例として説明する。したがって、本実施の形態において図1に示すコンピュータ装置は、プログラム制御されたCPU 101にて実現される通信制御手段及びネットワークインターフェイス106を介して、インターネットに接続し、ウェブサイトにアクセスする。

図2は、図1に示したコンピュータ装置にて実現される本実施の形態による情報検索システムの構成を示す図である。

図2に示すように、本実施の形態による情報検索システムは、インターネット上のウェブサイトからリンクを辿り所望の情報に関連のあるウェブサイトを検索する情報取得手段であるクローラ10と、クローラ10にて検索されたウェブサイトに対して所定の条件に基づく選別を行うウェブサイト選別部20と、ウェブサイト選別部20による選別にて選ばれたウェブサイトに基づいて、種々の戦略に基づくレポートを作成するレポート作成部30とを備えている。

#### 【0015】

上述したクローラ10、ウェブサイト選別部20及びレポート作成部30は、図1に示したメインメモリ103に展開されたプログラムにてCPU 101を制御することにより実現される仮想的なソフトウェアブロックである。CPU 101を制御してこれらの機能を実現させる当該プログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記憶媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供される。本実施の形態では、図1に示したネットワークインターフェイス106やフロッピーディスクドライブ108、図示しないCD-ROMドライブなどを介して当該プログラムを入力し、ハードディスク105に格納する。そして、ハードディスク105に格納されたプログ

ラムをメインメモリ 1 0 3 に読み込んで展開し、CPU 1 0 1 にて実行することにより、図 2 に示した各構成要素の機能を実現する。

#### 【 0 0 1 6 】

図 3 は、本実施の形態の情報検索システムによる情報検索の概略的な流れを示すフローチャートである。

図 3 に示すように、本実施の形態の情報検索システムは、図 2 に示したクローラ 1 0 により、初期サイトを獲得し（ステップ 3 0 1）、ユーザによって任意に選択された種々の戦略にしたがって動的にウェブサイトのクローリングを行う（ステップ 3 0 2）。

ここで、初期サイトとは、ウェブサイトのクローリングを開始するために初期的に設定されるウェブサイトもしくはその集合である。これをインデックスとして用いて次のクローリングを行う。また、戦略とは、ウェブサイトのクローリングを行う上での基準となる方針を意味し、具体的には検索条件等として設定される。本実施の形態で採られる戦略については 3 8 段落以下で詳述する。

次に、情報検索システムは、ウェブサイト選別部 2 0 により、クローリングで検索されたウェブサイトの集合の中から、検索条件であるトピックとの関連性や時間的条件に基づいて、有効なウェブサイトを選択する（ステップ 3 0 3）。そして最後に、ウェブサイト選別部 2 0 にて選択されたウェブサイトに対し、上記の戦略に基づいて評価を行い、レポートを作成する（ステップ 3 0 4）。作成されたレポートはウェブページ等の形でディスプレイ装置に表示され、あるいはハードディスク 1 0 5 等の記憶装置に保存される。

#### 【 0 0 1 7 】

本実施の形態において、クローラ 1 0 は、図 2 に示したように、初期サイトを取得する初期サイト獲得部 1 1 と、当該初期サイトに対応するウェブページに対して文書構造解析を行う文書構造解析部 1 2 と、文書構造解析部 1 2 による解析結果に基づいてクローリングによる取得対象であるウェブサイトの重要度を計算する重要度計算部 1 3 と、クローリングによるウェブサイトの取得処理を実行するクローリング実行部 1 4 とを備える。

初期サイト獲得部 1 1 における初期サイトの獲得には、例えば次のような方法

を採ることができる。

- ・情報を収集しようとする特定のウェブサイト（例えば企業サイト）のホームページ（トップページ）のURLを指定。
- ・任意のキーワードに対して既存の検索エンジンを用いて検索。

既存の検索エンジンとしては、インターネット上で提供される一般的な検索サービスを利用することができる。検索エンジンを用いて初期サイトを取得した場合は、検索によって取得されたウェブサイトの集合が初期サイト（初期サイト集合）となる。

#### 【 0 0 1 8 】

ウェブサイトのクロールは、初期サイトを起点として、そこからアンカータグなどの情報に基づいて参照されているウェブサイトを獲得し、その中から指定された戦略に一致するウェブサイトを獲得する処理である。この処理は、獲得されたウェブサイトの数の上限、深さの上限、クロールの実行時間などの終了条件を予め設定しておき、この終了条件を満たすまで再帰的に適用される。ユーザは、初期集合から動的にクロールを進めるためのヒントとして、個々のサイトとユーザが指定したトピックとの関連性を計算するために用いられるキーワードを指定する。クローラ 1 0 は、文書構造解析部 1 2 により初期サイト（初期サイト集合）に対応するウェブページの文書構造解析を行い、重要度計算部 1 3 によりクロールの際の戦略に応じたウェブサイト（HTML（Hypertext Markup Language）文書の記述上ではアンカータグ）の重要度を計算した上で、これらの情報を用いてクロール実行部 1 4 によりウェブサイトを検索し取得する。

#### 【 0 0 1 9 】

本実施の形態の文書構造解析部 1 2 による文書構造解析は、ウェブページのソースコードである HTML 文書においてブロックを識別する。

ここでブロックとは、特定の意味を持つ情報要素のまとまりであり、ウェブページを記述する HTML におけるブロックレベルとは必ずしも一致しない。このブロックに含まれる情報要素は、ブロックの属性（OBJECT\_LIST）に、情報要素のリストとして登録され、メインメモリ 1 0 3 や CPU 1 0 1 のキャッシュメモリに保存される。このブロック化により、HTML 文書に含まれる情報要素のう

ち、相互に関連する情報要素が対応付けられることとなる。

情報要素には、単一要素と複数の単一要素をマージして構成された複合要素とがある。解析の最初の時点では、全ての情報要素は単一要素として識別され、解析を進めることによって複数の情報要素が複合要素としてマージされる。この情報要素は、以下の属性を持つ。

- ・ TYPE : 情報要素のタイプ。

OBJECT\_ANCHOR : アンカー。

OBJECT\_TEXT\_BLOCK : テキスト。

OBJECT\_IMAGE : メディアのタイプ。他にAUDIO、VIDEOなどが定義可能。以降の説明では、全てのメディアタイプを代表してOBJECT\_IMAGEを記述する。

OBJECT\_DELIMITER : 情報要素のTYPEによらず区切り記号としての役割であると解析された場合に指定される。

- ・ URL : U R L。

OBJECT\_ANCHORの場合は、HREFで指定された値。

OBJECT\_IMAGE などの場合は、SRCなどで指定された値。

- ・ TITLE : タイトル。

OBJECT\_ANCHOR の場合は、A (アンカー) タグで囲まれたテキスト部分。

OBJECT\_IMAGE などの場合は、ALTなどで指定されたテキスト部分。

- ・ DESCRIPTION : 記述。

OBJECT\_TEXT\_BLOCK の場合は、そのテキスト部分。

OBJECT\_ANCHOR の場合は、関連するテキストが得られた場合、それらをマージする際にテキスト記述として追加される。

- ・ REFERRER : 参照情報。

OBJECT\_ANCHOR, OBJECT\_TEXT\_BLOCK の場合、それに関連する OBJECT\_IMAGE などの他のメディアタイプの情報要素が追加される。

- ・ EMPHASIS : 強調表現。

その情報要素が強調的な表現をされているかどうか指定される。

【 0 0 2 0 】

図 4 は、文書構造解析部 1 2 によるブロック化のアルゴリズムであるウェブペ



ージの文書構造解析処理の手順を説明するフローチャートである。

図 4 を参照すると、文書構造解析部 1 2 は、まず解析対象のウェブページの HTML 文書を読み込み、そのタグの記述に基づいて暫定的なブロック解析を行う（ステップ 4 0 1）。ここでは、HTML 文書において、一般的に意味が区切られる箇所に記述されることの多いタグを手掛かりにして、処理対象の HTML 文書を仮にブロックごとに区分する。このブロック解析に用いられるタグの例としては、例えば“BLOCKQUOTE”、“DD”、“DIV”、“DL”、“FORM”、“H1”、“H2”、“H3”、“H4”、“H5”、“HR”、“LI”、“P”、“TABLE”、“TD”、“TH”、“TR”、“UL”などが挙げられる。すなわち、これらのタグが出現した場合、ブロックの開始とする。そして、それぞれのタグに対応する終了タグ（</TABLE>など）が出現した場合は、そのブロックの終了とする。

また、上記のようにして決定されたブロック内に、次のようなタグ、あるいは、テキストが出現した場合は、そのブロック中に上述した情報要素を作成する。

A（アンカー）タグ：OBJECT\_ANCHORを生成する。URLとして、HREFで指定された値。TITLEとして、Aタグで囲まれたテキストを指定する。

タグで囲まれていないテキスト部分：OBJECT\_TEXT\_BLOCKを生成する。TITLEとして当該テキストを指定する。IMGなどメディアを指定するタグ：OBJECT\_IMAGEなどを生成する。IMGタグの場合、URLに、SRCで指定された値を、TITLEとして、ALTで指定された値を指定する。

さらに、これら個々の情報要素を作成する際に、その情報要素が強調表現であるかどうかを判断し、そうであれば、属性EMPHASISにその強さを指定する。強調表現であるかどうかは、例えば、“B”、“CENTER”、“EM”、“STRONG”、“TH”、“U”のようなタグによってその情報要素が囲まれているかどうかによって判断する。

以上のようにして得られたHTML文書の構造データは、図 1 に示したメインメモリ 1 0 3 や CPU 1 0 1 のキャッシュメモリに一時的に格納される。

【 0 0 2 1 】

次に、文書構造解析部 1 2 は、ステップ 4 0 1 で作成された構造データをメイ

ンメモリ 1 0 3 等から読み出し、この構造データに含まれる情報要素のうちで不要な情報要素を識別する（ステップ 4 0 2）。URLが同じであるOBJECT\_IMAGE、あるいは、同じTITLEを持つOBJECT\_ANCHOR、あるいは、同じDESCRIPTIONを持つOBJECT\_TEXT\_BLOCKが複数出現されている場合は、情報を運ぶという意味合いは弱く、単なる区切り記号的な役割を担う。したがって、不要な情報要素といえる。例えば、図 5（A）に示すようなイメージタグは、LI タグと同様の区切りの役割を果たしている。また、図 5（B）に示すようなアンカータグは、情報を運ぶ意味合いは少ない。

HTML 文書中に出現する情報要素の数を数えながら、このような情報を運ぶ意味合いの少ない情報要素を識別し、そのタイプをOBJECT\_DELIMITERとしてマークする。具体的には、不要イメージに関しては、以下のアルゴリズムによってマークを行う。なお、イメージの同一性は、URLが同じかどうかで判断する。

- ・ 1 つのブロックが 1 つのイメージのみを含む場合は、OBJECT\_DELIMITERの候補とする。
- ・ 1 つのブロックが複数の同一イメージのみを含む場合は、OBJECT\_DELIMITERの候補とする。
- ・ 文書全体において、複数のブロックで候補となっているものに対し、OBJECT\_DELIMITERをマークする。

また、不要なOBJECT\_ANCHOR、及びOBJECT\_TEXT\_BLOCKに関しては、同じタイトルを持つ要素が予め設定された所定の閾値以上あれば、それらをOBJECT\_DELIMITERとしてマークする。かかる閾値を用いた判断としては、例えば、テキスト長が、比較的短く（1 2 バイト以下）、かつ、出現回数が 3 回以上などであれば不要な要素とみなす。また、同じ文書内に非常に多く出現した（例えば 1 0 回以上）全く同じOBJECT\_ANCHOR、OBJECT\_TEXT\_BLOCKも不要な要素とみなす。これらの閾値は、システムの動作環境や用途などに応じて、経験的にあるいは適当な回数の実験を経て好ましい値を求めることができる。上記の例では、テキスト長は、実際のインターネット上のウェブページを対象として観察した結果から経験的に求めたものである（例えば、“先頭に戻る”、“戻る”、“キャッシュ”などがこの例に相当する）。また、出現回数の閾値も同様に実際に観察した結果から経験

的に求めたものである（例えば、購買サイトの“購入はこちら”などのアンカー）。

以上のようにして不要と判断されたアンカーは、クロールの対象とならない。

#### 【 0 0 2 2 】

次に、文書構造解析部 1 2 は、構造的に意味のないブロックを削除する（ステップ 4 0 3）。例えば、所定のブロックについて、そのブロック自身が情報要素を持たず、かつ内包するブロック（子ブロック）が 1 つの場合は、当該ブロックを下位ブロック（すなわち当該子ブロック）とマージする。

#### 【 0 0 2 3 】

次に、文書構造解析部 1 2 は、情報要素リストのマージ、ブロック分割及びリストタイプの識別を行う（ステップ 4 0 4）。まず、同一ブロック内の情報要素をマージする。情報要素のマージにより、関連する複数の要素が 1 つの複合要素とされる。以下、具体例を挙げる。

OBJECT\_ANCHORは、OBJECT\_ANCHOR、OBJECT\_TEXT\_BLOCK及びOBJECT\_IMAGEとのマージが可能であり、以下に示すような操作にてマージする。

- ・ OBJECT\_ANCHORとOBJECT\_ANCHORとのマージ：

2 つの情報要素のURLが同一の場合のみ行う。EMPHASIS属性が付与されている方を重要とみなし、そのTITLEを採用する。そして、EMPHASIS属性が付与されていない方は、そのTITLEをDESCRIPTIONに追加する。

- ・ OBJECT\_ANCHORとOBJECT\_TEXT\_BLOCKとのマージ：

OBJECT\_ANCHORのDESCRIPTIONにOBJECT\_TEXT\_BLOCKのDESCRIPTIONを追加する。

- ・ OBJECT\_ANCHORとOBJECT\_IMAGEとのマージ：

OBJECT\_ANCHORのREFERRERにOBJECT\_IMAGEを設定する。OBJECT\_ANCHORのDESCRIPTIONにOBJECT\_IMAGEのTITLEを追加する。

また、OBJECT\_TEXT\_BLOCKは、OBJECT\_ANCHOR、OBJECT\_TEXT\_BLOCK及びOBJECT\_IMAGEとのマージが可能である。なお、OBJECT\_ANCHORとのマージの場合の操作は上述したOBJECT\_ANCHORどうしの場合と同様であり、OBJECT\_IMAGEとのマージの場合の操作は上述したOBJECT\_ANCHORとOBJECT\_IMAGEとのマージの場合と同様で

ある。OBJECT\_TEXT\_BLOCKとOBJECT\_TEXT\_BLOCKとのマージの場合は、次に示すような操作にてマージする。

・ OBJECT\_TEXT\_BLOCKとOBJECT\_TEXT\_BLOCKとのマージ：

一方のOBJECT\_TEXT\_BLOCKのDESCRIPTIONに他方のOBJECT\_TEXT\_BLOCKのDESCRIPTIONを追加する。

基本的には、OBJECT\_ANCHOR、OBJECT\_TEXT\_BLOCK、OBJECT\_IMAGE が同一ブロック内に存在すれば、それらをまとめられるものと判断する。OBJECT\_ANCHORに関しては、同じURLを参照している要素が近くになれば、その中に含まれる要素をマージして1つの複合要素とする。

#### 【 0 0 2 4 】

また、同一ブロックに情報要素が3つ以下（OBJECT\_DELIMITERを含む場合は4つ以下）の場合、どのようにマージさせるかは容易に判断できるが、複数の情報要素がリストで並んでいるような場合は、情報要素リスト内の構造解析を行い、要素のマージ、あるいは、ブロックの分割を行う。情報要素リスト内の構造解析は、原則的にはNグラム統計を用いて次のように行う。

すなわち、1グラムから解析を行い、それぞれのグラム数において、同一ブロック内で支配的な要素の並びが見つかった場合、その並びで分割を行う。例えば、3グラム統計を行った場合、OBJECT\_DELIMITER、OBJECT\_ANCHOR、OBJECT\_TEXT\_BLOCKの当該ブロック図全体に対する割合（カバー率）が一定の閾値より大きい場合（例えば80%以上）は、その並びを分割し、情報要素をマージする。具体的には、図6に示す例では、2つの複合要素が作成される（bullet.gif（破線で囲んだ箇所）はOBJECT\_DELIMITERとして解析されている）。閾値を必ずしも100%としないのは、ブロック内に要素が列挙して書かれている場合でも、個々の要素に関しては、イメージがあったりなかったりするなどの揺れがあり、それを吸収するために経験的に閾値を決める必要があるためである（図6の例では、OBJECT\_DELIMITER、OBJECT\_ANCHOR、OBJECT\_TEXT\_BLOCKの並びは2回出現しており、それは全体をカバーしているので、カバー率はこの場合100%である）。

#### 【 0 0 2 5 】

ただし、Nグラム解析だけでは、解析できない場合がある。例えば、図7に示

す例のように、参考情報（破線で囲んだ箇所）が指定されている場合である。このような場合は、特にOBJECT\_ANHOR、OBJECT\_TEXT\_BLOCKに関する対応付けがNグラム統計では行うことができない。よって、情報要素に付与されているTITLE、DESCRIPTIONのテキスト部分からキーワードを抽出し、複数の情報要素間のキーワードリストの一致度（割合や数）を調べ、一致度が一定の閾値以上である場合は当該情報要素をマージする。例えば、双方のタイトルから抽出されたキーワードの一致する割合が、共に70%以上であるならば、マージする。一致度の計算の際にはキーワードの分類を考慮して、人名、組織名などの重みをより重くして一致度を計算することも可能である。例えば、ニュースなどの情報であれば特に人名、組織名が一致している場合その関連度は強い場合が多い。

#### 【0026】

また、図8に示す例では、複数の意味的なまとまりが1つのブロック内に存在する。破線で囲まれた各部分（“IBM関連リンク”、“日本IBM関連リンク”）がそれぞれ意味的なまとまりを構成している。このような場合は、タイトルとなる情報要素を区切りとして情報要素リストをブロックに分割する。タイトルとなる情報要素を識別するために、EMPHASIS属性を用いる。すなわち、一定の閾値以上の長さの情報要素リスト内にEMPHASIS属性が指定されているものが複数あり、その後の情報要素が類似した並びを持っている場合は、それを複数に分割する。この際、分割されたブロックのタイトルとして、EMPHASIS属性のタイトルを指定する。このような処理によって、図8のHTML文書に対して、図9に示すようなブロック分割がなされる。

一般に、情報要素リストの先頭の情報要素にEMPHASIS属性が付与されており、情報要素リスト中のOBJECT\_ANCHORが一定の閾値以上の割合である場合（例えば80%以上）は、当該EMPHASIS属性のタイトルをブロックのタイトルとして指定する。あるいは、先頭のみがOBJECT\_TEXT\_BLOCKであり、他の要素が一定の閾値以上（例えば80%以上）のOBJECT\_ANCHORである場合も同様に、先頭のOBJECT\_TEXT\_BLOCKのDESCRIPTIONをブロックのタイトルとする。これらの閾値は、システムの動作環境や用途などに応じて、経験的にあるいは適当な回数の実験を経て好ましい値を求めることができる。

## 【 0 0 2 7 】

このようにして抽出されたブロックのタイトルは、個々の情報要素に対する重要度の計算に用いられる。例えば図 8、9 に示した例で、トピックが“日本 I B M”であった場合、そのブロックに含まれる 3 つのアンカー“Home”、“プロダクト&サービス”、“サポート&ダウンロード”は、文字列“日本 I B M”を含まないが、ブロックのタイトルに含まれることによって関連する情報要素として判断される。このように文書構造を詳細かつ正確に解析することによって、単に位置的に近傍のテキストのみからでなく、離れた位置の要素間の依存関係を利用することが可能となる。

また、抽出されたブロックの情報要素リストにおいて、OBJECT\_ANCHORの割合が多い場合は、当該ブロックのリストタイプを次のように決める。

- ・ SITE\_MAP : リンク総数に対して、OBJECT\_ANCHORのホスト名が異なるリンクの数の割合が所定の閾値より小さい（例えば 5 0 % 以下）場合。それ以外の場合は、ホスト名に対して固有組織名を識別するための最小の文字列を求め、その文字列の異なるリンクの割合が所定の閾値より小さい場合（例えば 1 0 % 以下）は、SITE\_MAP とする。これらの閾値は、例えば、インターネット上の実際のウェブページに基づいて経験的に求められる。厳密にホスト名の一致を取るだけでは、たとえ同じ会社が提供しているページであっても異なる場合がある。ページの機能によって個別のホストを用いている場合があるからである。例えば、Yahoo! J apan(<http://www.yahoo.co.jp>) は、オークションには [auctions.yahoo.co.jp](http://auctions.yahoo.co.jp), 旅行関連は、[travel.yahoo.co.jp](http://travel.yahoo.co.jp) などホスト名を分けている。この場合は、固有組織を識別するための最小の文字列(yahoo.co.jp)の一致度によって、閾値に基づく判断を行うことができる。

- ・ LINK\_LIST : SITE\_MAPの条件を満たさない場合。

## 【 0 0 2 8 】

次に、文書構造解析部 1 2 は、以上のようにして得られたブロックごとの情報に基づいて、ブロック構造のマージを行う（ステップ 4 0 5）。

関連する情報要素は、必ずしもブロック内に連続して出現するわけではない。すなわち、ステップ 4 0 1 やステップ 4 0 4 の処理でブロックの設定を細かくし

過ぎてしまう場合がある。そこで例えば、所定のブロックの子ブロックにおいて並列構造があるならば、それを当該所定のブロック（親ブロック）の情報要素リストとすることにより、構造をマージする。マージする方法は、ステップ404で説明した情報要素リストに対するマージのアルゴリズムと同様である。

#### 【0029】

文書構造解析部12は、以上説明したステップ403乃至ステップ405の動作を、処理対象のHTML文書に対して適用が可能である限り繰り返して適用し、解析結果である構造データの構造が変更されなくなったならば、当該構造データをメインメモリ103やCPU101のキャッシュメモリに格納して構造解析の処理を終了する（ステップ406）。

以上のようにして、HTML文書の文書構造が意味のまとまりごとにブロック化され、HTML文書から抽出された情報要素がブロックの属性として記述されることにより、相互に関連する情報要素が対応付けられて当該ブロック内のアンカーに付加されることとなる。

#### 【0030】

次に、重要度計算部13によるクローリング先のウェブサイト（すなわちHTML文書におけるアンカータグ）の重要度の計算について説明する。

図10、11は、重要度計算部13により重要度を計算し、クローリング実行部14によりウェブサイトのコンテンツを取得する動作を説明するフローチャートである。

本実施の形態において、重要度計算部13の重要度の計算に用いられる基本的なアルゴリズム（図10、11にて説明されるアルゴリズム）は、下記文献2に開示されたFish-Searchと呼ばれる手法や、上述した文献1に開示されたShark-Searchと呼ばれる手法に基づく。

文献2：P. De Bra, G.-J Houben, Y. Kornatzki, and R. Post, Information retrieval in distributed hypertexts, in Proceedings of RIAO'94, Intelligent Multimedia, information Retrieval systems and managements, New York, NY, 1994.

ただし、本実施の形態による手法は、ユーザにより指定された戦略及び文書構造解析部 1 2 により解析されてアンカーに付加された情報要素に基づいて検索対象であるウェブサイトのスコア（重要度）を計算し、このスコアに応じてクロールの対象を動的に決定する点、及びユーザに対してキーワードリストを提示することによって、よりユーザの目的にあったウェブサイトを収集するためのインタラクションを可能にした点が拡張されている。

また、クロールの候補となるのは、文書構造解析部 1 2 による解析の結果として抽出された、意味のある複合情報要素のうち、他の文書への参照を持つもののみである。

#### 【 0 0 3 1 】

図 1 0、1 1 を参照すると、まず、ユーザが本実施の形態による情報検索システムを構成するコンピュータ装置の入力手段を操作することにより、パラメータの設定及び初期設定が行われる（ステップ 1 0 0 1）。具体的には、初期ノード集合（初期サイト、以下、クロールするウェブサイトをノードと称す）、探索幅（width）、探索の深さ（D、depth）、初期ノード集合のサイズ（S）、時間制限、探索用のキーワード（Domain Query、Focused Query）、戦略（STSET）などを設定する。戦略の設定においては複数の戦略を選択することが可能であり、各戦略に対して重み付けを行うことができる。また、クローリングの回数として 0（crawlingCount = 0;）がセットされる。

次に、重要度計算部 1 3 が、初期ノード集合の個々のノードの深さをパラメータ D にセットし、それらを空のリスト（以下、ノードリスト）に挿入する（ステップ 1 0 0 2）。また、メインメモリ 1 0 3 等から文書構造解析部 1 2 による解析結果である構造データを読み出す。そして、ノードリストが空でなく、処理されたノードの数がパラメータ S より小さく、かつ時間制限内である間、次の処理を繰り返し実行する（ステップ 1 0 0 3）。

#### 【 0 0 3 2 】

まず、クローリングの回数を 1 加算（crawlingCount += 1;）する（ステップ 1 0 0 4）。ここで、crawlingCount が一定の増分を越えた場合（例えば 1 0 0



サイトごと)、かつstrategicScoreの計算が大域的である場合に、strategicScoreの再計算を行い、ノードリスト中のスコア(個々のノードのスコア、すなわち各ノードであるウェブサイトへのリンクするアンカーの重要度)の値を置き換える(ステップ1005)。また、後述するように関係するキーワードが抽出されるものに関しては、それらのキーワードリストを提示してユーザによる選択を促す。ユーザがキーワードの選択を行った場合は、選択されたキーワードに応じてトピックを更新する。この間、クロールは一旦中断しても構わないが、続けることも可能である。

次に、ノードリストから先頭のノードを取り出し、カレントノードとする。このカレントノードは、図1に示したメインメモリ103やCPU101のキャッシュメモリに保持され、クロール実行部14にて読み出される。そして、クロール実行部14がネットワークインターフェイス106を介してインターネットにアクセスし、当該カレントノードのURLを持つコンテンツ(ウェブページや種々のデータ等)を取得する(ステップ1006)。取得したコンテンツは、図1に示したメインメモリ103やハードディスク105などの記憶装置に格納される。

#### 【0033】

次に、重要度計算部13は、カレントノードに関する探索の深さ(パラメータD)を調べ、If depth > 0ならば、以下の手順でカレントノードの関連性を計算する(ステップ1007)。

まず、カレントノードからリンクされているノード(以下、子ノードと称す)のスコア(child\_node.inherited\_score)を計算する(図11、ステップ1008)。この計算手順は以下の通りである。

```
If relevance(current_node) > 0
Then child_node.inherited_score = d * strategicScoreForPage(STSET, current_node)
```

dは、予め定義された減衰定数。0より大きく1より小さい。

```
Else child_node.inherited_score = d * current_node.inherited_score
```

次に、`child_node.anchor_score`を計算する（ステップ1009）。

`child_node.anchor_score = (relevance(anchor) + strategicScoreForAnchor(STSET, anchor))/2`

そして、子ノードのpotential scoreを計算する（ステップ1010）。

`child_node.potential_score = g * child_node.inherited_score + (1 - g) * child_node.anchor_score`

`g`は、予め定義された定数。0より大きく1より小さい。

【0034】

次に、重要度計算部13は、カレントノードの全ての子ノードに対して以下の計算を行う（ステップ1011～1013）。

If 子ノードが優先リスト中に存在しているか。

Then

i) そのノードに対するリスト中の値と今計算されたpotential\_scoreの大きい方を求める。

i i) スコアを最大値で置き換える。

i i i) 子ノードをリスト中の適切な位置に移動する。

Else `child_node`にpotential\_scoreを付け、リスト中の適切な位置（スコア順）に挿入する。

【0035】

さらに、重要度計算部13は、カレントノードの全ての子ノードに対して以下の計算を行う（ステップ1014）。

深さ（`child_node.depth`）を計算する。

If カレントノードが関連がある場合

Then `child_node.depth = D`

Else `child_node.depth = current_node.depth - 1`

If 子ノードが優先リスト中に存在している。

Then

そのノードに対するリスト中の値と今計算された深さの大きい方を求める。

その値で置き換える。

## 【 0 0 3 6 】

ステップ 1 0 1 4 までの処理が終了した後、またはステップ 1 0 0 7 で、If  $\text{depth} > 0$  でない場合は、ステップ 1 0 0 3 に戻り、各条件を満足する限り、クローリングの回数を 1 加算してステップ 1 0 0 4 以降の処理を繰り返す。そして、ステップ 1 0 0 3 のいずれかの条件を満足しない場合は、重要度計算部 1 3 及びクローリング実行部 1 4 による処理を終了する。

## 【 0 0 3 7 】

次に、上述したアルゴリズムにおける個々の計算方法について説明する。

・  $\text{relevance}(\text{current\_node})$  の計算方法

Domain Query と Focused Query は、ベクトル（トピックベクトル）で表現される。そして、これらの Query（キーワード）とテキストの一致度は、ベクトル間の距離（内積など）で計算される。 $\text{current\_node}$  は、そのテキスト部分をベクトルに変換し、その類似度を計算する。ユーザの指定によって、関連がないと判断されたキーワードは、トピックベクトルにおいてマイナスの重要度を持つ。これらは、以下の式で計算される。

$$\text{relevance}(\text{current\_node}) = \text{Similarity}(\text{current\_node}, \text{Domain Query}) + \text{Similarity}(\text{current\_node}, \text{Focused Query})$$

・  $\text{relevance}(\text{anchor})$  の計算方法

$$\text{relevance}(\text{current\_node}) = \text{Similarity}(\text{TITLE}, \text{Domain Query}) + \text{Similarity}(\text{DESCRIPTION}, \text{Focused Query})$$

・  $\text{strategicScoreForPage}(\text{STSET}, \text{current\_node})$  の計算方法

ユーザによって指定された戦略ごとのスコアの重み付き総和によって決定される。値は 0 から 1 の間に正規化する。個々の戦略に対するスコアの計算方法（ $\text{strategicScoreForPage}(\text{ST}, \text{current\_node})$ ）は後述する。

・  $\text{strategicScoreForAnchor}(\text{STSET}, \text{anchor})$  の計算方法

ユーザによって指定された戦略ごとのスコアの重み付き総和によって決定され

る。値は 0 から 1 の間に正規化する。個々の戦略に対するスコアの計算方法 (`strategicScoreForPage(ST, current_node)`) は後述する。

【 0 0 3 8 】

次に、本実施の形態によるクローリングで用いられる戦略とそのタイプ、及び `strategicScoreForPage(ST, current_node)`, `strategicScoreForAnchor(ST, anchor)` の計算方法について、例を挙げて説明する。

本実施の形態で用いられる戦略には、局所的なものと大局的なものの 2 つのタイプがある。局所的な戦略は、ウェブページ内の情報のみで重要度を決定できるが、大局的な戦略は、複数のウェブページを解析することによって重要度を計算する。

- ・ ユーザが指定したトピックに近いウェブサイトを検索（局所的戦略）

この戦略は、`relevance(current_node)`, `relevance(anchor)` によって計算されるものであり、図 1 0、1 1 に示した基本的なアルゴリズム内に組み込まれている。

【 0 0 3 9 】

- ・ 重要なウェブサイトを検索（大局的戦略）

これは、多くのウェブサイトで同じ情報が提供されていれば、その情報は重要だと見なす戦略である。同じ情報があるかどうかを複数のサイト内で調べる必要があるので大局的である。複数のウェブサイトで同じ情報が提供されているかどうかは、例えば次の文献 3 に開示されているような、ウェブページからヘッドライン（見出し）を抽出する公知技術を用いることによって、知ることができる。

文献 3：武田、野美山、”サイト・アウトライニング-インターネットからの情報収集と可視化技術-”，情報処理，Vol. 42，No.8，2001.

この方法では、結果として同じ事柄を言及した情報要素の集合のみを返す。例えば、図 1 2 に示すようなヘッドラインを持つウェブサイトが検索されることとなる。なお、図 1 2 において、Site2～4はSite 1の子ノードである。またこの方法では、抽出されたヘッドラインを構成する情報要素のテキスト部分からキーと

なるキーワードとその重みを抽出して特徴キーワードリストを生成する。図 1 2 の例では、主要な要素となるキーワード“ロータス”、“チボリ”、“日本IBM”、“統合”とその重みとが特徴リストにリストアップされる。

ヘッドラインに含まれる新しい情報要素に対しては、対応ノードのテキスト部分から抽出されたキーワードリストと、ヘッドラインの特徴キーワードリストの距離（内積などで求められる）を計算し、その距離を重要度とする。

この方法において、`strategicScoreForPage(ST, current_node)`は、カレントノードのウェブページに含まれるヘッドラインの重要度の総数に基づいて計算される。そして当該ウェブページのアンカーの総数で割ることによって正規化する。

また、`strategicScoreForAnchor(ST, anchor)`は、現在抽出されているヘッドライン集合の全ての特徴ベクトルとの距離の最大値とする。

#### 【 0 0 4 0 】

・関連するイメージなどを多く含むウェブサイトを検索（局所的戦略）

HTML文書の文書構造を解析することによって、テキストと関連する他のメディアタイプ（Image、Audio、VideoあるいはMIME（Multipurpose Internet Mail Extensions）タイプに定義されているドキュメントファイル（例えばPDF（Portable Document Format））など）を対応付けることが可能となる。テキスト部分がトピックと関連性を持つかどうかに基づいて重要度が計算される。

この方法において、`strategicScoreForPage(ST, current_node)`は、カレントノードのウェブページに含まれる関連イメージの総数であり、当該ウェブページのイメージの総数で割ることによって正規化する。

また、`strategicScoreForAnchor(ST, anchor)`は、関連するイメージであれば、重要度を1とし、そうでない場合は重要度を0とする。

#### 【 0 0 4 1 】

・重要なキーワードを含む情報に基づく検索（大局的戦略）

テキストの情報から抽出されたキーワードに基づいてクラスタリングを行い、各クラスターで重要と判断されたキーワードを多く含むかどうかによって重要度を判断する。この方法については、特開 2 0 0 1 - 3 2 5 2 7 2 公報に詳細に開示

されている。

この方法において、`strategicScoreForPage(ST,current_node)`は、カレントノードのウェブページに含まれるホットワードを含む要素の重要度の総数に基づいて計算される。そして当該ウェブページのアンカーの総数で割ることによって正規化する。

また、`strategicScoreForAnchor(ST,anchor)`は、所定の要素がホットワードを含むのであれば、その重要度（0以上1以下）の値に設定し、そうでない場合は0とする。

#### 【 0 0 4 2 】

##### ・アンカーの数に基づく検索（局所的戦略）

カレントノードのウェブページ内に存在するアンカーの総数によって重要度を判断する。

この方法において、`strategicScoreForPage(ST,current_node)`は、カレントノードのウェブページに含まれるアンカーの総数に基づいて計算される。局所的に重要度を計算するには、例えば、リンクの数を11段階に分け（0、1、2、3、4、5、6、7、8、9以上）、それぞれにスコア（0から1までの値）を与える。大局的に計算する場合は、検索されたウェブサイトの集合において、最大数のアンカーを含むウェブページの当該アンカーの数で正規化する。

また、`strategicScoreForAnchor(ST,anchor)`は、全て0とする。

#### 【 0 0 4 3 】

##### ・期間の限定に基づく検索（局所的戦略）

ウェブページ内に出現する情報の発信された期間を限定する。HTTPプロトコルで得られる最終更新日付が、限定された期限内であれば、重要度を1とし、そうでなければ、期限からの超過日数を正規化した値を重要度とする。

##### ・リンク集に基づく検索（局所的戦略）

HTML文書の文書構造解析によって得られたリンク集に含まれるリンクの数を重要度とする。

##### ・被参照リンクの数に基づく検索（大局的戦略）

他のウェブサイトから参照されている数を重要度とする。

・被参照リンクの参照の数（大局的戦略）

カレントノードのウェブページを参照しているリンク数を重要度とする。

【 0 0 4 4 】

上記のようにして、クローラ 1 0 は、ユーザにより選択された任意の戦略を用いてウェブサイトのクローリングを行い、検索されたウェブサイトの集合（以下、サイト集合）を得る。得られたサイト集合は、ステップ 1 0 0 6 で個々に取得されたコンテンツと共に、メインメモリ 1 0 3 やハードディスク 1 0 5 等の記憶装置に格納される。

【 0 0 4 5 】

ウェブサイト選別部 2 0 は、メインメモリ 1 0 3 等の記憶装置に格納されているサイト集合の中から不要なウェブサイトを選別し、選別されたウェブサイト及びそのコンテンツを削除する。不要なウェブサイトとしては、トピックとの関連がないウェブサイト、時間条件を満たさないウェブサイトが挙げられる。

クローラ 1 0 はユーザにて指定されたトピックを表すキーワードに基づき所定の戦略に従ってクローリングを行うのであるが、得られたサイト集合には、ユーザが指定したトピックと関連が低い、あるいは関連のないウェブサイトが含まれる可能性がある。そこで、クローリングにより得られたサイト集合に対してトピックとの一致度を求め、関連のないウェブサイトはサイト集合から削除する。ただし、クローリング時に得られる参照構造で関連すると判定された複数のウェブサイトの間に位置するウェブサイト（当該ウェブサイトを介してリンクされる場合）は、中間位置の当該ウェブサイト自体にトピックとの関連性がない場合でも削除しない。

また、サイト集合には、ユーザが指定した時間に一致しないウェブサイト（指定期限内に検索されていないウェブサイトなど）が含まれる可能性がある。そこで、クローリングにより得られたサイト集合に対して HTTP プロトコルで得られる LAST\_MODIFIED 属性とユーザが指定した時間条件とを比較し、一致しないウェブサイトはサイト集合から削除する。

【 0 0 4 6 】

レポート作成部 3 0 は、クローラ 1 0 にて検索され、ウェブサイト選別部 2 0

にて選別されたサイト集合に対して、個々の戦略に対するスコアの総得点を計算し、それらを正規化したレポートを作成し、メインメモリ 1 0 3 やハードディスク 1 0 5 等の記憶装置に格納する。例えば、イメージに関しては、1 サイト当たりの関連イメージの数などが情報としてレポートに含まれることとなる。レポートは、例えば HTML 文書として作成し、ウェブブラウザを用いて閲覧できるようにすることができる。

複数のトピックに対する戦略のスコアと比較する（平均値との比較、あるいは標準偏差を求める）ことにより、そのトピックがどのような戦略に一致しているのか（例えば、イメージが多い、情報が多い（掲示板などが多く含まれる））などの傾向を知ることができる。

#### 【 0 0 4 7 】

以上のようにして、本実施の形態の情報検索システムによれば、最適な戦略を組み合わせることによって、ユーザの目的により適したサイト集合を獲得することができる。

また、クローリングするために用いられるアンカーが独立でなく、対応するテキストに対して適切に対応付けられるため、トピックとの関連性をより正確に判断することができる。

さらに、文書構造解析を行って、HTML 文書内で意味のないアンカーを排除するため、不要なウェブサイトをクロールする無駄を防ぐことができる。

そして、文書構造解析によって得られたブロックを利用することにより、位置的に離れた要素の依存関係をクローリングに利用できる。

このブロックを認識することによって、ウェブページ中のリンク集を特定することができるため、高品質なリンク集を収集し抽出することもできる。

この他、クロール中に（例えば図 1 0 のステップ 1 0 0 5 の段階で）ユーザに対して関連キーワードを提示することによって、指定されたトピックの曖昧性を解消できる。例えば、トピック「ジャガー」は、「車」「動物」「ミック・ジャガー（ロック歌手）」などの曖昧性を持つが、クラスタリングなどによって、関連キーワードが表示し、ユーザがこの関連キーワードを指定することにより、クロール対象を絞り込むことができる。具体的には、車の「ジャガー」を検索した



い場合は、「ミック・ジャガー」、「ライブ」などの関連キーワードを指定してマイナスの重要度を与えることによって、これらの関連キーワードを含むウェブサイトをクロール対象から外すことができ、結果としてトピックの曖昧性を解消することができる。

#### 【 0 0 4 8 】

図 1 3、1 4、1 5 は、具体的なウェブページの HTML 文書に対して文書構造解析を行った様子を示す図である。図 1 3 はウェブページを表示した様子を示し、図 1 4 は、図 1 3 のウェブページの HTML 文書に対して図 4 のステップ 4 0 1 におけるタグによるブロック解析を行った状態の構造データを示し、図 1 5 は、ステップ 4 0 2 以降の解析処理により整理されたブロック列の状態の構造データを示す。なお、図 1 4、1 5 には HTML 文書の解析結果の一部のみを記載している。

図 1 3 を参照すると、対象である HTML 文書は、レイアウトを揃えるためにテーブル (TABLE) タグを多用していることがわかる。そのため、これらのオブジェクト 1 4 0 1 ~ 1 4 0 6 は、表示上はまとまって見えるとしても内部構造上は離れている場合がある。例えば、画面の下部中央にある「ニュース」オブジェクトは、図 1 3 に示すように表示上はまとまっているが、実際の構造は図 1 4 に示すようにテーブルタグを用いて位置揃えがなされており、意味的には余分なタグが多数挿入されている。本実施の形態の文書構造解析を行うことにより、図 1 5 に示すように、これらのオブジェクトが 1 つの要素 1 5 0 1 として解析され、さらにそのタイトルとして「ニュース」が付与される。

図 1 4 における dotted\_rule\_197px.gif のイメージオブジェクト 1 4 0 7 は、区切りとして利用されているが、これも正しく認識されており、解析中は区切りの意味合いを持つ情報として利用されるが、図 1 5 に示す解析後の構造には含まれない。

また、e-business hosting ( HYPERLINK "http://www.ibm.com/services/jp/webhosting/" http://www.ibm.com/services/jp/webhosting/) のオブジェクト 1 4 0 8 ~ 1 4 1 0 などのように、同じ URL を指している情報要素は、本実施の形態の文書構造解析により、図 1 5 に示すように 1 つの要素 1 5 0 2 にまとめら

れる。

【 0 0 4 9 】

このように、文書構造解析部 1 2 による解析の結果、HTML 文書の構造が意味のあるまとまりごとにブロック化されるため、この構造に基づいてアンカーのリンク先へ遷移し、ウェブサイトクロールすることにより、無駄のない適切なクロールを行うことができる。

【 0 0 5 0 】

【発明の効果】

以上説明したように、本発明によれば、情報の使用目的に応じて多様な戦略による柔軟な情報検索を可能とすることができる。

また、本発明によれば、ウェブページのクロールにおいて、この多様な戦略による情報検索を実現するために、ウェブページに含まれる情報を有効に活用して検索を行うことができる。

【図面の簡単な説明】

【図 1】 本実施の形態による情報検索システムを実現するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

【図 2】 図 1 に示したコンピュータ装置にて実現される本実施の形態による情報検索システムの構成を示す図である。

【図 3】 本実施の形態の情報検索システムによる情報検索の概略的な流れを示すフローチャートである。

【図 4】 本実施の形態の文書構造解析部によるウェブページの文書構造解析処理の手順を説明するフローチャートである。

【図 5】 本実施の形態の文書構造解析処理により解析される不要な情報要素の例を示す図である。

【図 6】 本実施の形態の文書構造解析処理による情報要素のマージの例を説明する図である。

【図 7】 本実施の形態の文書構造解析処理による情報要素のマージの他の例を説明する図である。

【図 8】 本実施の形態の文書構造解析処理による情報要素のマージのさら

に他の例を説明する図であり、マージを行う前の状態を示す図である。

【図 9】 図 8 の例における情報要素のマージを行った状態を説明する図である。

【図 1 0】 本実施の形態の重要度計算部により重要度を計算し、クローリング実行部によりウェブサイトのコンテンツを取得する動作を説明するフローチャートである。

【図 1 1】 本実施の形態の重要度計算部により重要度を計算し、クローリング実行部によりウェブサイトのコンテンツを取得する動作を説明するフローチャートである。

【図 1 2】 同じ事柄を言及したウェブサイトのサイト集合の例を示す図である。

【図 1 3】 ブラウザにて表示されたウェブページの例を示す図である。

【図 1 4】 図 1 3 のウェブページの HTML 文書をタグによりブロック解析した状態の構造データを示す図である。

【図 1 5】 図 1 4 の状態からさらに文書構造解析を行った状態の構造データを示す図である。

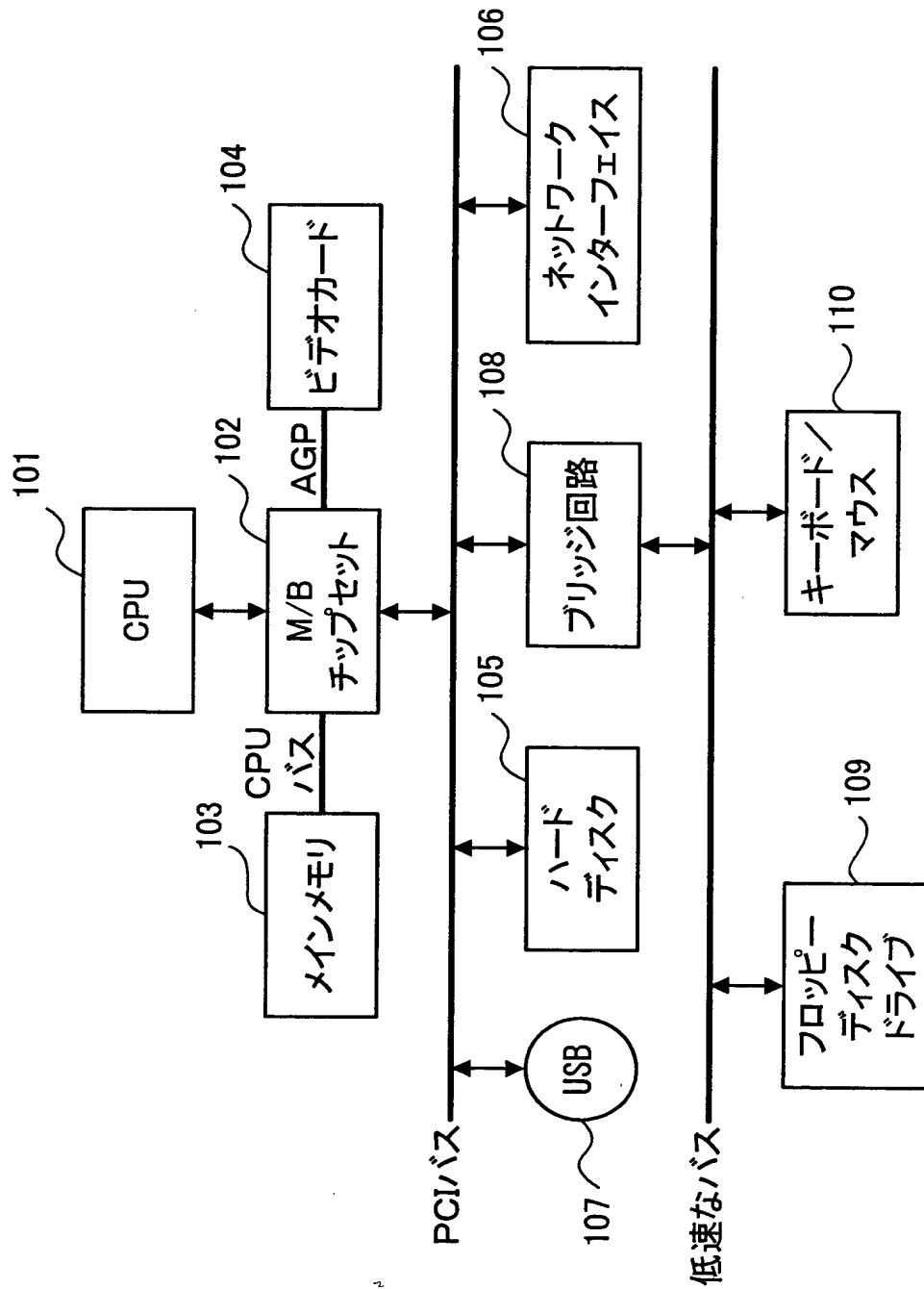
【符号の説明】

1 0 … クローラ、1 1 … 初期サイト獲得部、1 2 … 文書構造解析部、1 3 … 重要度計算部、1 4 … クローリング実行部、2 0 … ウェブサイト選別部、3 0 … レポート作成部、1 0 1 … CPU、1 0 2 … M/B チップセット、1 0 3 … メインメモリ、1 0 5 … ハードディスク、1 0 6 … ネットワークインターフェイス

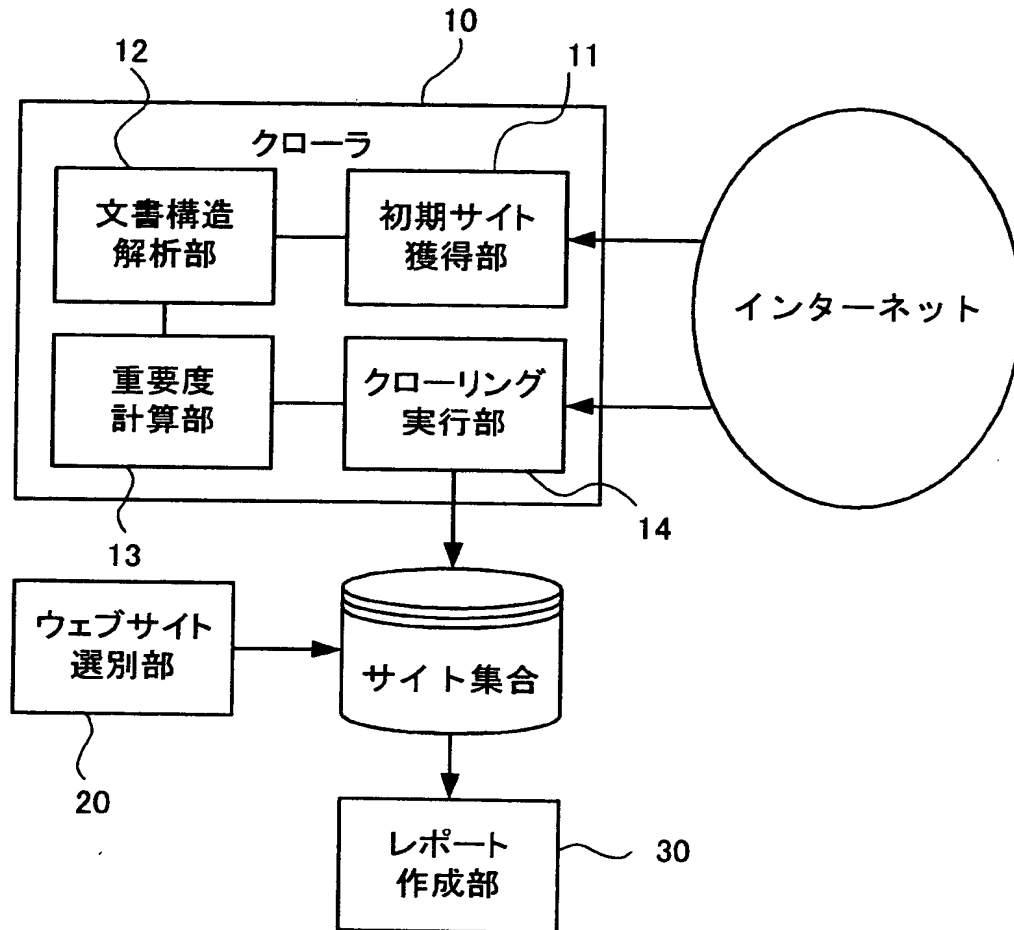
【書類名】

図面

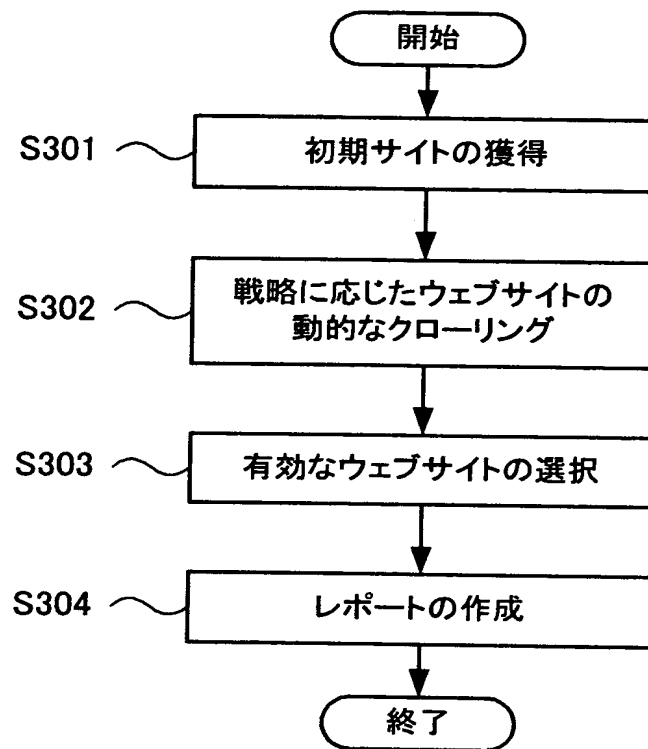
【図 1】



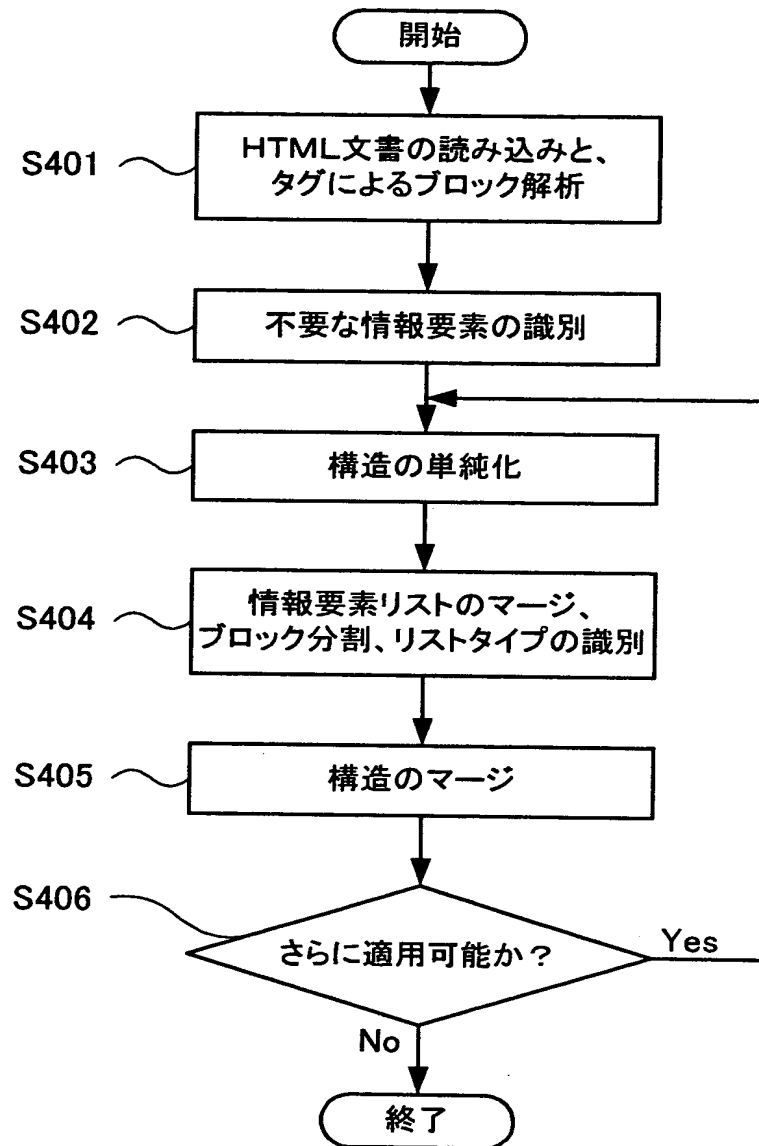
【図 2】



【図 3】



【図 4】



【図 5】

(A)

```
⋮  
項目 1<br>  
項目 2<br>  
⋮
```

(B)

```
⋮  
<a href="top.htm">先頭に戻る</a>  
⋮  
<a href="top.htm">先頭に戻る</a>
```



【図 6】

```
|
<a href="news1.htm">developerWorks最新情報</a>
今週の新コンテンツは、Web services、Java technology、Linux、Open
source、XMLの各ゾーンに追加されました。
|
<a href="news2.htm">インターネット翻訳の王様バイリンガル Version5
Linux版、デリバリー(ダウンロード販売)を開始</a>
文脈により適切・的確な翻訳を実現する「インテリジェント翻訳エンジン」
搭載した「インターネット翻訳の王様バイリンガル Version5」。
```

【図 7】

```

<a href="news1.htm">developerWorks 最新情報</a>
今週の新コンテンツは、Web services、Java technology、Linux、Open source、XML の各
ゾーンに追加されました。
「詳細は、<a href="detail.htm">こちら</a>を参照してください。」

<a href="news2.htm">インターネット翻訳の王様バイリンガル Version5 Linux版、デリバ
リー(ダウンロード販売)を開始</a>
文脈により適切・的確な翻訳を実現する「インテリジェント翻訳エンジン」搭載した「インター
ネット翻訳の王様バイリンガル Version5」。
```

【図 8】

```

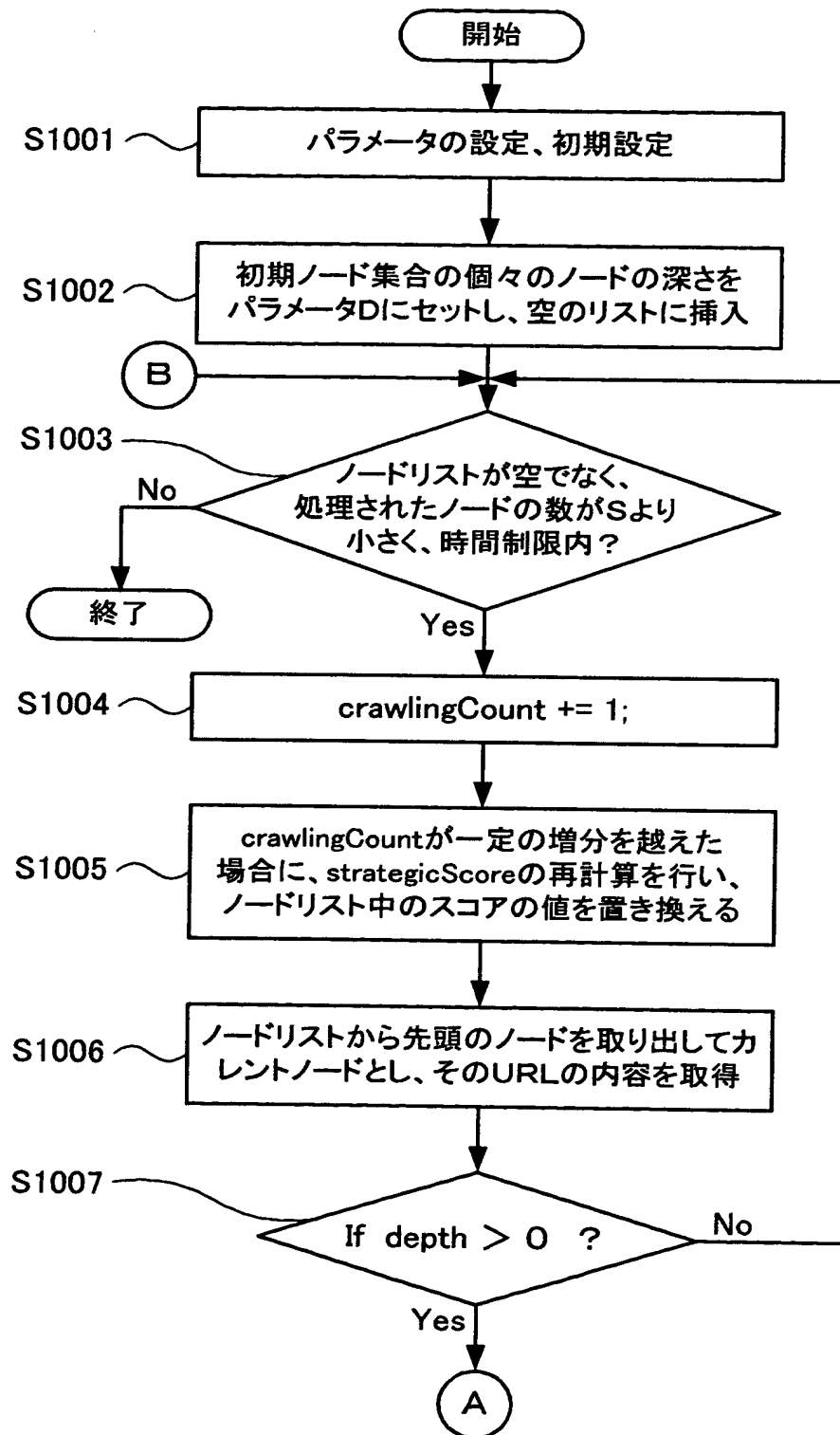
<b>IBM 関連リンク(US)</b>
<a href="http://www.ibm.com">Home</a><br>
<a href="http://www.ibm.com/products/us">Products & Services</a><br>
<a href="http://www.ibm.com/support">Support & downloads</a><br>

<b>日本IBM 関連リンク(US)</b>
<a href="http://www.ibm.com/jp">ホーム</a><br>
<a href="http://www.ibm.com/products/jp">製品 & サービス</a><br>
<a href="http://www.ibm.com/jp/support">サポート & ダウンロード</a><br>
    
```

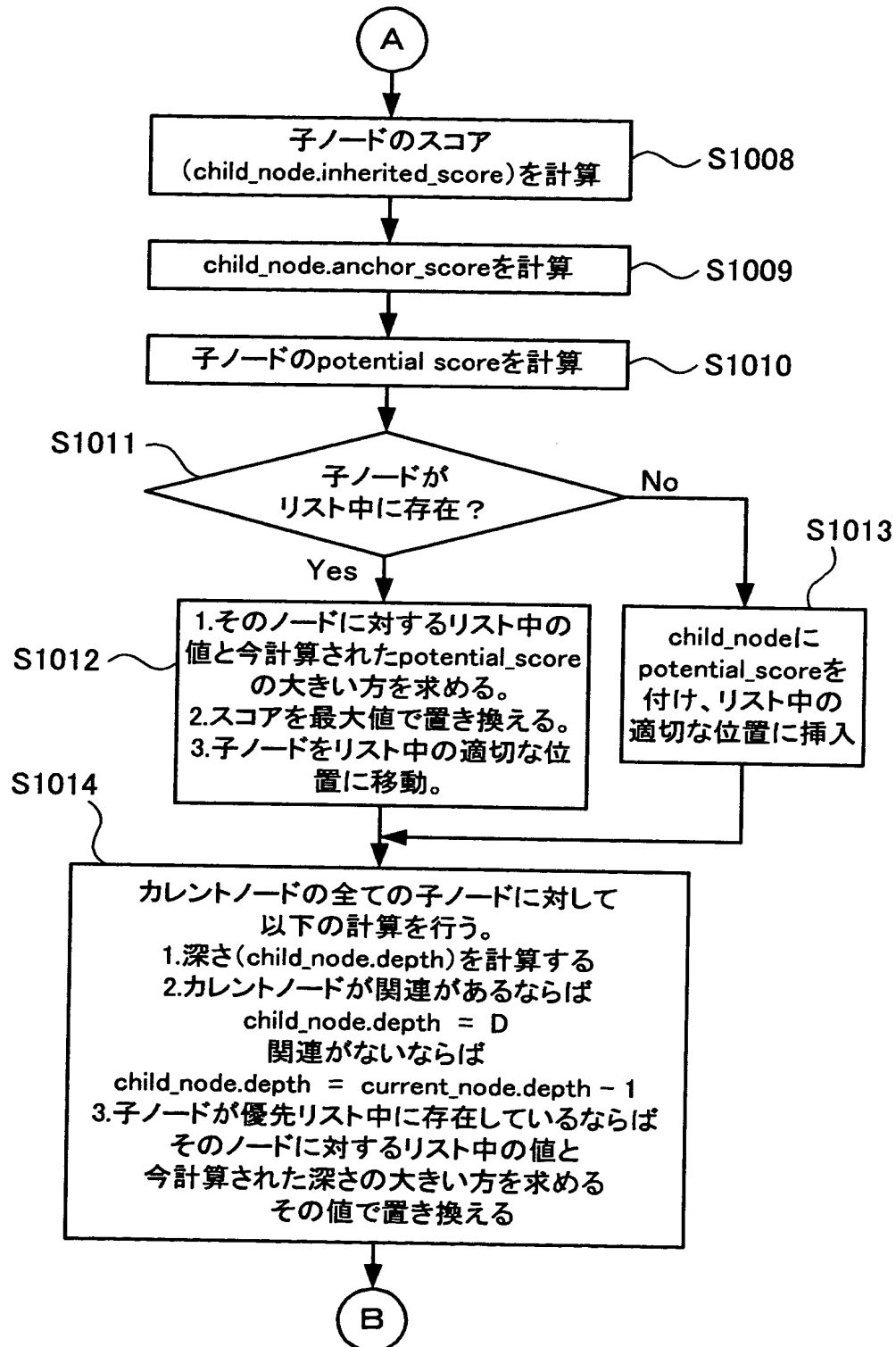
【図 9】

```
<block>
<block title="IBM関連リンク">
  <object type="anchor" url="http://www.ibm.com" title="Home"/>
  <object type="anchor" url="http://www.ibm.com/products/us" title="Products & Services"/>
  <object type="anchor" url="http://www.ibm.com/support" title="Support & downloads"/>
</block>
<block title="日本IBM関連リンク">
  <object type="anchor" url="http://www.ibm.com/jp" title="Home"/>
  <object type="anchor" url="http://www.ibm.com/products/jp" title="プロダクト & サービス"/>
  <object type="anchor" url="http://www.ibm.com/jp/support" title="サポート & ダウンロード"/>
</block>
</block>
```

【図 1 0】



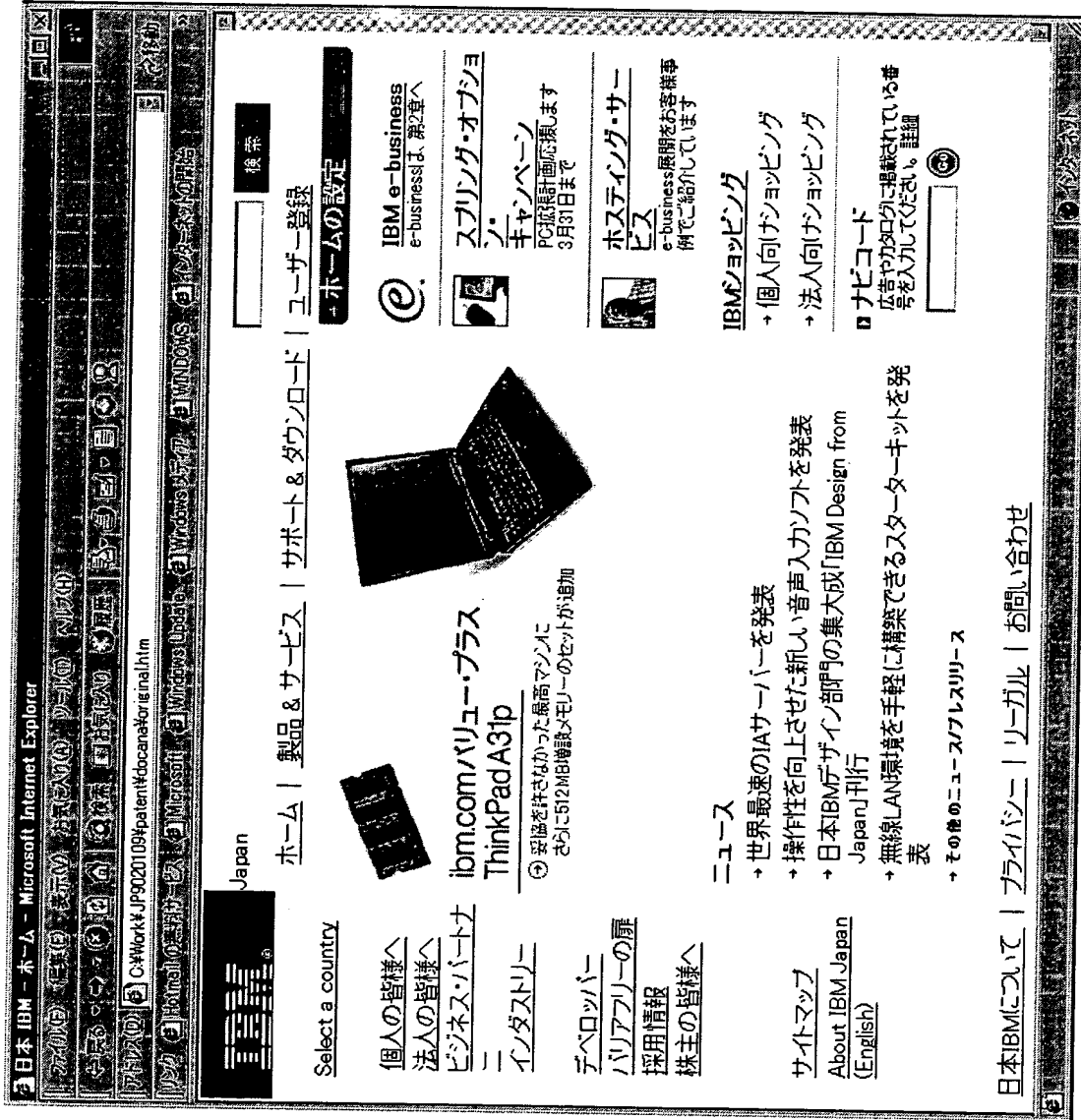
【図 1 1】



【図12】

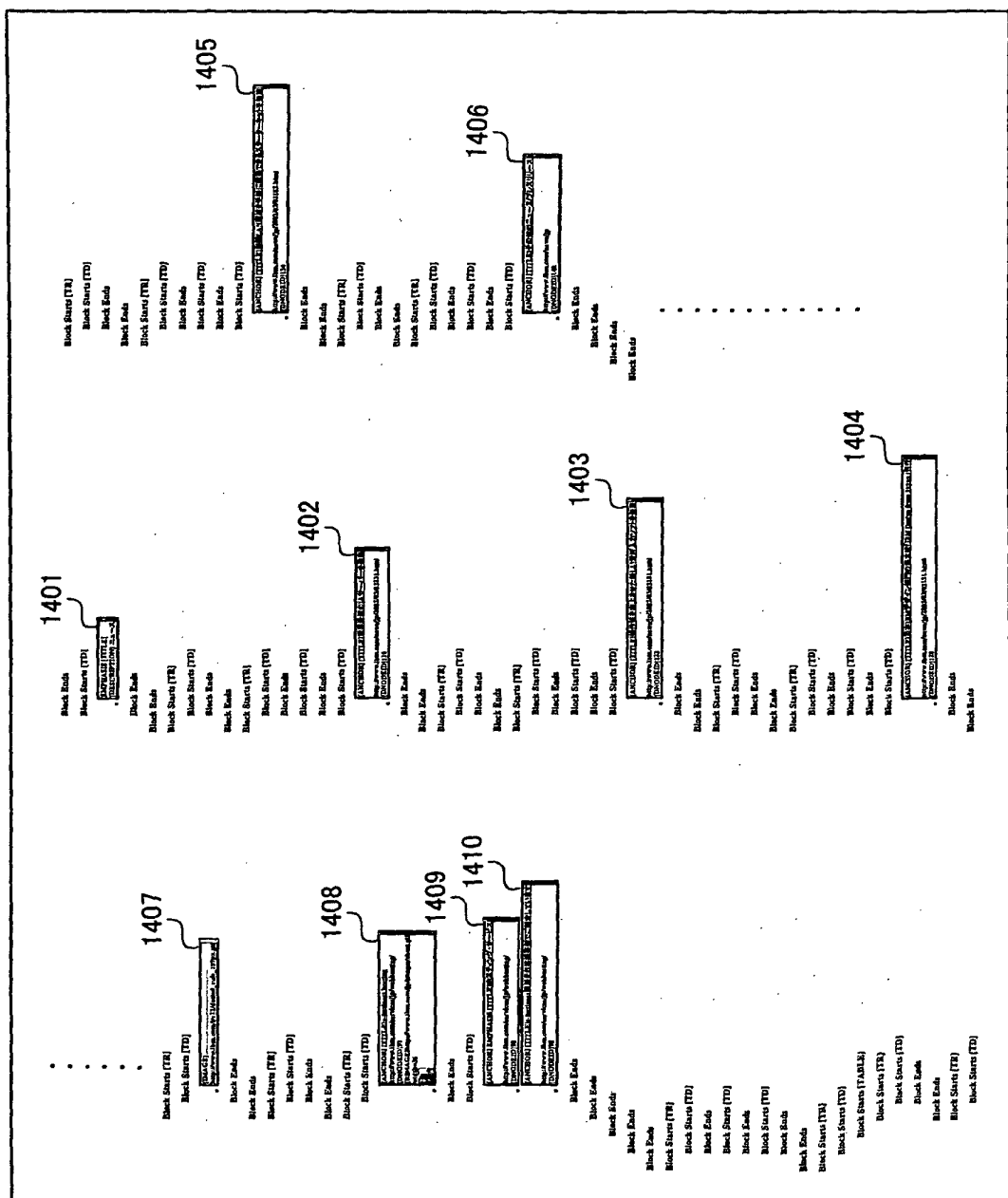
この夏、ロータスと日本チボリステムズが日本IBMに統合 [Site 1]  
+日本IBM、ロータスらを7月に統合 [Site 2]  
+日本IBM、ロータス、チボリを統合 [Site 3]  
+ロータスとチボリ、日本IBMに統合 [Site 4]

【図 13】

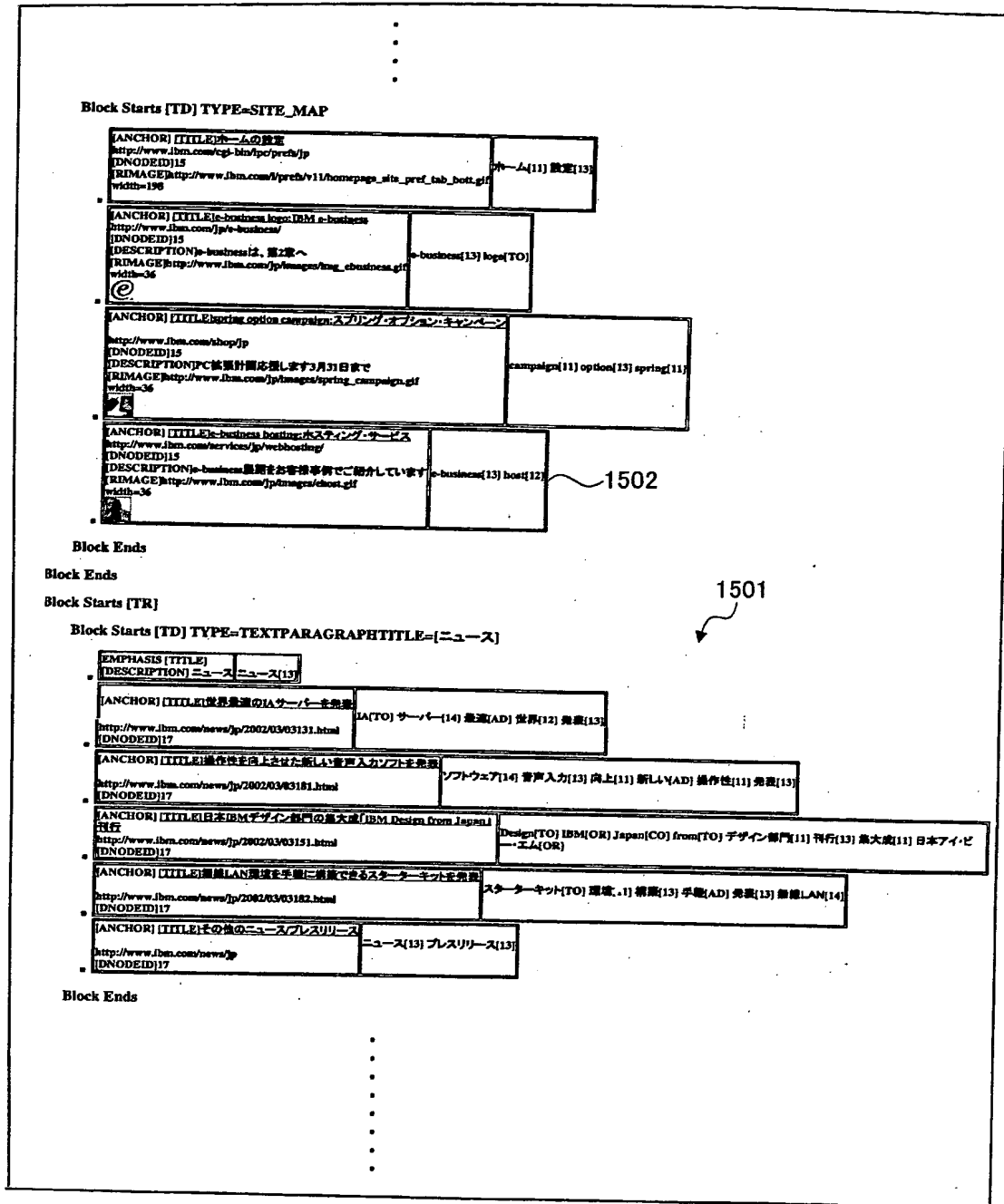




【图 14】



【図15】



【書類名】                      要約書

【要約】

【課題】    コンピュータを用いた情報検索において、情報の使用目的に応じて多様な戦略による柔軟な情報検索を効果的に実現する。

【解決手段】    所定のウェブページにおける意味を考慮してHTML文書の構造を解析する文書構造解析部12と、この解析結果に基づき、予め定められた戦略にしたがって、このウェブページからリンクされる他のウェブサイトの重要度を計算する重要度計算部13と、この重要度計算部13により計算された重要度に応じてウェブサイトをクロールするクロール実行部14とを備える。

【選択図】                      図2

## 認定・付加情報

|         |                          |
|---------|--------------------------|
| 特許出願の番号 | 特願 2 0 0 2 - 2 1 1 6 3 4 |
| 受付番号    | 5 0 2 0 1 0 6 6 7 4 9    |
| 書類名     | 特許願                      |
| 担当官     | 土井 恵子 4 2 6 4            |
| 作成日     | 平成 1 4 年 1 0 月 8 日       |

### <認定情報・付加情報>

#### 【特許出願人】

|          |  |
|----------|--|
| 【識別番号】   | 390009531                                      |
| 【住所又は居所】 | アメリカ合衆国 1 0 5 0 4、ニューヨーク州 アーモンク ニュー オーチャード ロード |
| 【氏名又は名称】 | インターナショナル・ビジネス・マシーンズ・コーポレーション                  |

#### 【代理人】

|          |   |
|----------|---|
| 【識別番号】   | 100086243                                       |
| 【住所又は居所】 | 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内 |
| 【氏名又は名称】 | 坂口 博  |

#### 【代理人】

|          |   |
|----------|---|
| 【識別番号】   | 100091568                                       |
| 【住所又は居所】 | 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内 |
| 【氏名又は名称】 | 市位 嘉宏   |

#### 【代理人】

|          |   |
|----------|---|
| 【識別番号】   | 100108501                                     |
| 【住所又は居所】 | 神奈川県大和市下鶴間 1 6 2 3 番 1 4 日本アイ・ビー・エム株式会社 知的所有権 |

|          |       |
|----------|-------|
| 【氏名又は名称】 | 上野 剛史 |
|----------|-------|

#### 【復代理人】

申請人

|          |   |
|----------|---|
| 【識別番号】   | 100104880                                     |
| 【住所又は居所】 | 東京都港区赤坂 5 - 4 - 1 1 山口建設第 2 ビル 6 F セリオ国際特許事務所 |

|          |       |
|----------|-------|
| 【氏名又は名称】 | 古部 次郎 |
|----------|-------|

出 願 人 履 歴 情 報

識別番号 [390009531]

1. 変更年月日 2002年 6月 3日

[変更理由] 住所変更

住 所 アメリカ合衆国10504、ニューヨーク州 アーモンク ニ  
ュー オーチャード ロード

氏 名 インターナショナル・ビジネス・マシーンズ・コーポレーショ  
ン